

杭州亚运会网络舆情的主题挖掘和传播特征研究

——基于LDA主题模型

周海燕, 周璐

(湖北大学 体育学院, 湖北 武汉 430062)

摘要: 目的和意义: 探究亚运会舆情的发展走势, 识别并确定出亚运会期间的热点主题, 从而了解赛事舆情的传播特征和情感演变过程, 为引导体育赛事网络舆情的治理提供思考。方法和过程: 以信息生命周期理论为理论基础, 以第十九届杭州亚运会体育赛事为研究对象, 借助Python爬虫软件获得齐全的博客数据, 通过词频分析得到有效的语料, 同时进行LDA主题建模以确定主题最优数目并进行主题挖掘, 探究体育赛事网络舆情的传播特征。结果和结论: 体育赛事网络舆情的用户地域分布具有明显差异性, 体育赛事网络舆情的微博评论文本相较于传统文本更具有娱乐化特征, 体育赛事网络舆情易出现群体的情感极化现象, 体育赛事网络舆情能表达族群认同和身份认同感。结合体育赛事网络舆情自身存在的问题, 相关部门应创新舆情引导形式, 转变舆情管理角度, 提高网络舆情预判能力。

关键词: LDA主题模型; 杭州亚运会; 网络舆情; 传播特征

Research on Theme Mining and Communication Characteristics of Network Public Opinion of Hangzhou Asian Games: Based on LDA Topic Model

ZHOU Hai-yan, ZHOU Lu

(School of Physical Education, Hubei University, Wuhan 430062, China)

Abstract: Objective and Significance: To explore the development trend of the public opinion of the Asian Games, identify and determine the hot topics during the Asian Games, so as to understand the communication characteristics and emotional evolution process of the public opinion of the event, and provide governance thoughts to guide the network public opinion of sports events. Method and Process: Based on the information life cycle theory and the 19th Hangzhou Asian Games sports event as the research object, complete Weibo data were obtained by using python crawler software, effective corpus was obtained by word frequency analysis, and LDA theme modeling was carried out for optimal number of topics and topic mining, so as to explore the communication characteristics of online public opinion of sports events. Result and Conclusion: The geographical distribution of users of online public opinion of sports events is significantly different. Weibo comment text of online public opinion of sports events is more entertaining than traditional text. Online public opinion of sports events is prone to the phenomenon of group emotional polarization, and online public opinion of sports events expresses ethnic identity and identity. Combined with the problems of online public opinion of sports events, relevant departments should innovate the form of public opinion guidance, change the perspective of public opinion management, and improve the ability to predict online public opinion.

Key words: LDA topic model; Hangzhou Asian Games; Online public opinion; Propagation characteristics

习近平总书记在党的二十大报告中就舆论工作发表重要讲话:加强全媒体传播体系建设,推动形成良好网络生态^[1]。中国互联网络信息中心(CNNIC)在京发布第52次《中国互联网络发展状况统计报告》显示,截至2023年6月,我国网民规模达10.79亿人,较2022年12月增长1109万人,互联网普及率达76.4%^[2]。随着互联网和人工智能的快速发展,我国早已进入自媒体时代,各大社交媒体平台亦成为大众价值观念交锋和抒发情绪的主阵地。社交媒体平台的交互性和便捷性,打破了传统媒体传播速度受限的壁垒,大众可以迅速获取信息、表达意见,这为社会热点事件在网络上的传播和舆情发酵提供了较为宽松和广阔的空间,使得舆情事件影响更为广泛。在某些情况下,关键意见领袖的观点甚至会左右整个舆情事件的发展方向。体育赛事是备受全民关注的热点事件,话题拓展性极强,容易成为众多媒体和大众关注的焦点。国家体育总局提出:切实肩负起体育媒体舆论主阵地职责使命,为杭州亚运会、巴黎奥运会备战,为体育强国建设营造良好舆论氛围^[3]。大型体育赛事的网络舆情能够塑造一个国家的形象和声誉,影响体育治理体系和治理能力,也会影响运动员的口碑和职业生涯。中国作为杭州亚运会的东道主,展现了国家的良好形象和组织能力,增加了亚运会在国内外的曝光度和讨论热度,网络上关于亚运会的各种话题和争议层出不穷,形成了广泛的舆论场,因此有必要全面研究杭州亚运会网络舆情,防范和应对未来体育赛事网络舆情的潜在风险。

本文基于文本挖掘技术,构建LDA主题模型,挖掘杭州亚运会各个阶段舆情的发展走势,识别主题,从而了解赛事舆情的传播特征和情感演变过程。相关研究结果能够为未来体育赛事网络舆情走向及治理提供参考依据,为营造健康的体育网络环境提供支持。

1 研究方法与设计

1.1 LDA主题挖掘模型

LDA(Latent Dirichlet Allocation)是机器学习和自然语言处理等领域常用的文本挖掘方法,主要用于从文本中发现隐含的、有用的信息的方法^[4]。它通过无监督学习的方式,不需要对数据人工标注,就可以对文本主题及其分布进行可视化分析。相较于传统的文本分析方法,其更具客观性和主题区分度。

LDA主题模型是由文档、主题、词项(特征词)三个层次构成的三层贝叶斯模型,形成“文档—主题”和“主题—词”两个重要的概率分布。“文档—主题”概率矩阵能够揭示每篇文档中不同主题所占的比例,“主题—词”分布则可进行主题识别和凝练,根据每个主题下的特征

词来判断主题的内容。

1.2 样本选取

杭州亚运会是亚洲顶级体育盛事,中国作为主办国,旨在展示国家形象。本文以亚运会为例,分析体育赛事网络舆情,预测其趋势,并利用新浪微博数据研究网络舆情传播,目的是为未来大型赛事网络舆情管理提供建议,减少负面影响。

1.3 数据获取

本文运用Python爬虫软件抓取微博评论文本。选定2023年9月23日至10月9日为时间范围,通过微博高级检索功能,抓取了#杭州亚运会#及相关热搜话题下的微博数据,共计191,000条评论。数据包括用户信息、发布时间、内容及互动量等。由于评论数量大且存在无效信息,接下来将对这些数据进行预处理。

1.4 文本预处理

对文本进行预处理是决定整个实验数据质量好坏的关键一环,会影响后续数据挖掘的准确性和完整性。将初步爬取的杭州亚运会数据集导入Python并存储,对数据进行清洗,具体步骤如下:(1)去重降噪处理。去掉表情emoji、标点符号以及特殊字符,去掉带超话或话题标签,去掉用户的重复评论,去掉评论中无效的数字和英文字母,保留文本中有效的汉字、数字和英文字母。

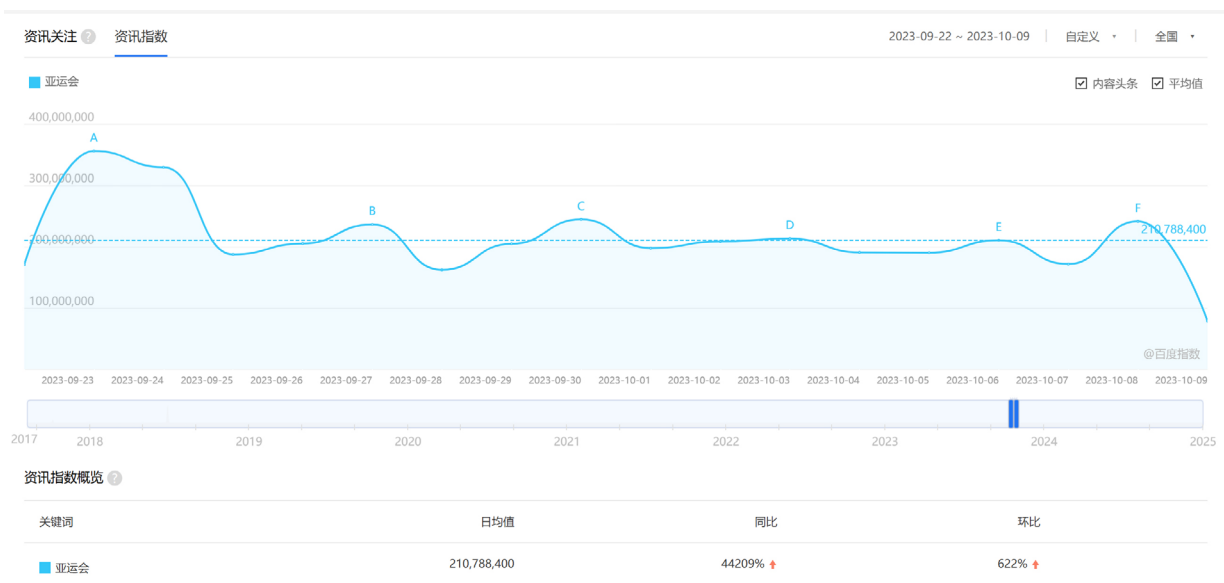
(2)文本分词。采用Jieba库中的精确模式对微博用户的转发、评论信息进行切分,形成有实际意义的词语。

(3)去停用词。本文使用哈工大停用词表对文本进行加工过滤,过滤无关字符,去除不影响微博用户表达主题倾向的无效字词。(4)进行词频统计。

1.5 网络舆情时间演化走势

本文依据百度指数数据,分析了杭州亚运会的网络关注度和媒体报道趋势。搜索指数反映了用户对亚运会的关注和支持变化,资讯指数则体现了新闻报道的关注度和持续性。根据网络舆情周期理论,以天为单位,杭州亚运会的网络生命周期为2023年9月23日至10月8日。

在亚运会预热期,网民关注度逐渐上升,开幕式时达到高潮,彰显了亚运会的盛大规模以及观众对开幕式的高度认可。首金产生后的14天正式比赛期间,舆情讨论出现波动,表明这一阶段更受关注,符合体育赛事媒体关注的周期性变化。闭幕式时,相关词条登上热搜,主创团队推出的“亚运花园”“数控草坪”和“数字火炬人”引发了大量讨论和热度。总体来看,杭州亚运会的网络舆情走势呈现出“两头高,中间低且有起伏”的特点,与以往大型体育赛事的报道规律有所不同。



(数据来源: 百度指数)

图1 杭州亚运会搜索指数



(数据来源: 百度指数)

图2 杭州亚运会资讯指数

1.6 网络舆情词频分析

词频分析法是一种常用的文本挖掘方法。利用词频分析可以描述和预测产业、事物的发展趋势,判断事物之间的关联性^[5]。分析文本数据中的高频词汇有助于揭示公众关注点以及对热点事件的态度。这种分析不仅应用于图书情报研究,也适用于网络舆情分析。作为中国主要的社交媒体平台,微博对公众认知有着重大影响。通过分析微博上关于杭州亚运会的热搜话题和评论,可以发现体育赛事相关的网络舆情焦点。在2248条评论中,

有301个不同词汇出现,其中前100个高频词的分布已进行可视化展示。

图3展示了由Wordcloud2生成的词云图,该图通过可视化手段呈现了微博评论文本中词语的出现频率。在图中,出现频率较高的词语以较大的字体显示,其中词频超过1000次的词汇依次为“亚运会”“杭州”“中国”和“决赛”,具体出现次数分别为41401次、37523次、13273次和7746次。这些高频词汇在一定程度上揭示了杭州亚运会的主要关注点,也反映了公众对亚运会选手的积极支持态度。

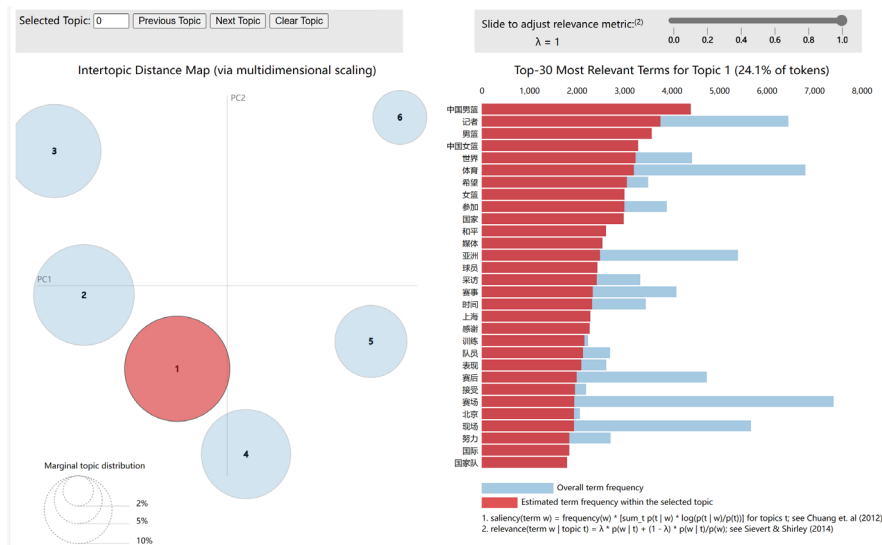


图5 基于 pyLDAvis 的主题可视化气泡图

有研究表明,当 λ 接近1时,该主题下呈现的30个关键词与主题呈正相关^[8]。图5中左侧的气泡分布分别代表6个不同的主题,各主题之间的气泡距离体现了主题间的差异度,气泡有重叠说明2个主题的特征值有交叉^[9]。研究显示六个主题气泡互不重叠,表明主题识别效果良好且各具特色。气泡大小体现主题在语料库中的占比,主题1最大,主题6最小。每个主题旁列有30个关键词。选择主题时,面板展示与之最相关的30个关键词,红色显示所选主题关键词频次,蓝色显示全部关键词频次。主题气泡直径显示标记概率,从高到低依次为:Topic1(24.1%),Topic2(22%),Topic3(18.9%),Topic4(17.3%),Topic5(11.3%),Topic6(6.3%)。

3) LDA 主题模型训练

完成“杭州亚运会”LDA模型分析后,使用gensim库对预处理的微博文本进行训练。选取各主题前20高频词,显示主题区分效果良好。分析得出6个主题及其特征词:主题1关注我国体育优势以及亚洲对手表现;主题2关注王一博等名人对亚运会传播的影响;主题3比较中国男篮、女篮表现;主题4关注金牌争夺;主题5关注网友对开闭幕式的期待;主题6关注肖战助力亚运。(如表1所示)

表1 杭州亚运会网络舆情“主题—特征词”分布

主题	特征词
Topic 1	韩国 王楚钦 樊振东 晋级 半决赛 日本 乒乓球 女团 战胜 英雄 男团 男单 联盟 女单 羽毛球 国乒 教练 组合 印度 击败
Topic 2	王一博 霹雳舞 宣传片 正当燃 全红婵 电竞 王者 荣耀 滑板 推广 郑钦文 街舞 领衔 能量 期待 传递 热烈 首次 打卡
Topic 3	中国男篮 记者 中国女篮 世界 体育 希望 参加 国家 和平 媒体 亚洲 球员 采访 赛事 上海 感谢 训练 表现 赛后 接受

续表

主题	特征词
Topic 4	孙颖莎 吴艳妮 夺冠 成绩 100 夺金 游泳 祝贺 铜牌 张雨霏 结束 银牌 奖牌 接力 恭喜 本届 200 首金 陈雨菲 田径
Topic 5	开幕式 新闻 赛场 中国 女排 健儿 央视 现场 网友 浙江 期待 闭幕式 火炬 亚洲 数字 观众 精彩 人民日报 代表团
Topic 6	肖战 喜欢 品牌 少年 中国 女足 梦想 未来 期待 祝福 圆梦 大使 演员 青年 全球 平安 担任 青春 圆满成功

2 讨论与分析

2.1 基于 LDA 模型的亚运会网络舆情主题内容分析

运用 LDA 模型对微博评论进行分析,根据每个主题的前20个高频词,回顾并分析文本数据中的多义词问题。结合杭州亚运会相关微博话题、帖子和评论,对 LDA 模型所识别的评论主题进行深入内容分析。

主题1涉及韩国、王楚钦、樊振东、日本、乒乓球和羽毛球等关键词,展现了我国在这些传统优势项目上的竞争态势。韩国和日本作为亚洲的主要竞争对手,近年来竞技水平显著提升,由此激发了网民的民族自豪感和对中国运动员的期待。国民对于此类大型体育赛事的关注度和讨论热度颇高,这反映出网民对高水平体育比赛的情绪回应。王楚钦和樊振东作为颇具影响力的国际运动员,其表现会直接影响舆情走向,而任何失误都可能被放大。当他们表现不佳或比赛失利时,网民的情绪反应尤为强烈且非理性。

主题2突出了王一博、霹雳舞、全红婵、王者荣耀等关键词。王一博作为中国滑板运动大使演唱了火炬传递主题曲《燃》,并担纲霹雳舞宣传片的领衔舞者,引发了全网热议。其影响力帮助街舞项目吸引了更多目光,

扩大了该项目的影响力。同时,杭州亚运会首次将电子竞技列为正式比赛项目,吸引了全球电竞爱好者的关注,提升了亚运会的社会影响力。电竞作为体育赛事的国际化进程得到了加速,尤其是具有国家文化特色的场地体验,给人们带来了强烈的视觉冲击。据报告,有3.5亿人观看了亚运电竞比赛,相关讨论量达到400万条,中国在电子竞技项目上取得了四金一铜的佳绩,赢得了网民的广泛赞誉。此外,全红婵作为跳水界的焦点人物,其在亚运会上的表现也引发了广泛关注和网络舆情,人们期待她能为中国队赢得更多荣誉。

主题3涉及中国男篮和女篮在杭州亚运会的表现。男篮在半决赛中最后一节表现不佳,以一分之差输给菲律宾,未能进入决赛,引发球迷不满。相比之下,女篮在与日本的比赛中表现出色,成功卫冕冠军。两队的表现差异引发了网友的广泛关注和讨论,相关话题如#亚运中国三大球不尽如人意#、#男篮心态失衡#、#中国女篮冲金#等成为网络热点。

主题4涉及孙颖莎、吴艳妮、夺冠、成绩和祝贺等关键词。在女子双打输给印度后,网络上涌现负面情绪,但孙颖莎在国庆节战胜日本选手,赢得赞誉。与此同时,在女子100米栏决赛中,吴艳妮因抢跑被取消成绩,引发热议。尽管有负面评论,但许多网友鼓励吴艳妮继续努力。体育赛事结果之所以引发关注,与比赛的竞争性紧密相关,体现了体育的独特魅力。观众对其所支持运动员的胜负特别关注,期望他们能够赢得比赛。

主题5可以看到开幕式、赛场、中国女排、数字、期待、闭幕式等关键词,杭州亚运会的开闭幕式赢得了网友们的广泛喜爱和期待,特别是开幕式上3D动画技术营造的数字烟花以及闭幕式中的“亚运花园”“数字草坪”等元素,不仅展现了中国传统文化的魅力,也增强了国家形象和国际影响力。此外,“中国女排”成为热议的焦点,显示出公众对女排姑娘们的高度期待和对比赛的极大关注,相关话题频繁出现在热搜榜上。女排之所以备受网民青睐,不仅是因为她们的出色表现,更因为她们所代表的“中国女排精神”——坚韧不拔、勇敢奋斗、无所畏惧,这一精神已成为中国共产党人精神谱系中极具体育特色的重要组成部分。在大型体育赛事中,“女排”始终是网民热议的核心话题之一。

主题6可以看到肖战、喜欢、品牌、中国女足、梦想、公益等关键词,可以凝练出肖战担任杭州亚运会圆梦公益大使这一话题。经笔者对比,发现主题6与主题2都有共现词“期待”,出现1838次。体育与娱乐的结合日益紧密,明星为体育赛事带来流量,响应国家政策,推动“互联网+全民健身”新趋势,创新方式以扩大体育受众。中国女

足虽在半决赛中不敌日本,但展现了坚韧斗志,虽未晋级决赛,却荣获铜牌,赢得广泛赞誉。足球作为全球最受欢迎的运动,中国女足承载了球迷对国足的殷切期望。

2.2 杭州亚运会网络舆情的传播特征分析

1) 用户地域分布不均

在亚运会网络舆情研究过程中,通过百度指数分析微博用户地域信息时发现,华东和华南地区是主要的舆情发源地,尤其是杭州、北京、上海等经济文化发达地区。微博文本词频统计显示,浙江、杭州、北京、上海、广东等地地名频繁出现,这与地区发展程度密切相关。同时,东部沿海城市的明星运动员如苏炳添、张雨霏、樊振东等,因家乡情怀而成为讨论焦点。

东部沿海地区因经济和文化水平较高,其网络舆论传播热度高于其他地区。在大型体育赛事中,舆情管控重点针对这些地区,尤其是人口密集、复杂度高的城市^[10]。同时,应考虑各城市的文化特点,例如杭州作为亚运会主办城市,其网络信息阅览和传播行为习惯应作为舆情管控的首要考量因素。

2) 评论文本娱乐化

新闻娱乐化是媒体顺应市场需求的表现,导致了新闻内容和表达方式的改变,使之更贴近大众生活,挖掘娱乐元素,进而增强受众的接受度和趣味性。网络信息技术精彩地转播各种赛事,并将运动员最新的采访和动态传达给观众,同时在网络化、多媒体化背景下的竞技体育新闻更多地呈现出娱乐化、话题化倾向。在整个网络社会被娱乐思潮席卷以及体育新闻自身娱乐化的双重影响下,体育赛事网络舆情同样呈现出娱乐化的表达方式^[11]。观看体育比赛已成为人们娱乐生活的一部分,微博评论的娱乐化表达正影响着体育话语风格。社交媒体平台的包容性使得微博评论语言更为有趣生动,体育爱好者常用幽默语言和网络流行语来表达对运动员的支持以及对赛事的正面情感。例如,杭州亚运会微博中出现的“显眼包”“yyds”等词汇,彰显了网友对运动员的称赞和亲切感。夸张的表达方式,如将亚运会背景音乐称为“适合中国宝宝的BGM”,也体现了网友对赛事的情感共鸣,促进了信息交互和舆情形成,丰富了赛事内容,提升了亚运会的影响力。

3) 群体的情感极化

网络的去个人化特质有利于群体认同的形成,在消弭群内个人差异的同时,也扩大了群间的差异,使得网络舆论容易朝着原有偏向的方向持续发展,最终呈现出极端的观点^[12],这就是舆论极化。互联网上并非人人都能时刻保持清醒与理性,一旦他人的言论契合自身某种认知框架而产生强烈的情感,就会在传播过程中循环

往复,从而引发群体情感极化^[13]。

微博为用户提供了表达和对话的空间,但也有不理性的行为出现。例如,在亚运会上,中国队因汪雪儿犯规被取消成绩,部分微博用户发表负面言论攻击运动员。吴艳妮的抢跑事件也引发了网络暴力,支持和反对的群体形成对立观点,导致情感极化。即便如此,杭州亚运会整体网络舆情仍以正面情绪为主。因此,应促进个体与群体间的情感认同,避免情感极化。

4) 民族认同感

杭州亚运会微博评论中,正向情感词汇如“加油”和“祝贺”占据主导地位,同时地域性名词如“中国队”和“日本”频繁出现。体育赛事作为全球文化的重要组成部分,能让不同国家的观众在不同地点共享赛事,激发民族情感。亚运会期间,社交媒体上人们表达族群观点,形成情感共鸣。乒乓球女单决赛中日对决,中国队获胜时,微博话题热度达到顶峰,体现了群体团结和民族认同,强化了民族自豪感和归属感。

3 研究结论与建议

依据杭州亚运会网络舆情的传播特征等研究结果,本文针对有关部门进行体育赛事网络舆情监管和引导提出如下建议。

(1) 转变应对网络舆情分析角度

相关部门需从传统媒体时代的网络舆情分析角度转变到自媒体时代的网络舆情角度,摒弃媒体对舆论完全可控这一落后的观念意识,在体育赛事网络舆情的传播方式和传播效果等方面做出有效改变。

(2) 加强数智技术手段应用

研发智能监管系统。借助自动化、智能化手段对整个体育比赛期间的舆情进行监管和研判,提升监管效率和准确性。针对体育事件相关的社交媒体、新闻发布、公众评论等舆论信息进行实时监测和预警,及时察觉并处理可能引发社会不稳定的负面舆情信息。

(3) 提升舆情引导能力

加强正面宣传。弘扬正能量,控制赛事负面情绪的传播,提高公众对体育新闻的辨识能力,引导公众营造正确的网络舆论氛围。此外,扩大意见领袖在社群中的影响力,借助他们的言论和行为来引导网络舆情,形成积极向上的舆论环境。

(4) 增强多元主体协同治理能力

鼓励公众参与体育事件的网络舆情监管和引导工作,拓宽公众参与渠道,可通过设立微博举报平台、开展问卷调查等方式,收集公众对体育话题的意见和建议。同时,加强舆论监督,发挥媒体、公众等多主体协同治理的作用,

对体育赛事网络舆情监管和引导工作进行监督和评价,推动相关部门不断改进和完善工作。

参考文献

- [1] 新华网. 习近平:高举中国特色社会主义伟大旗帜为全面建设社会主义现代化国家而团结奋斗:在中国共产党第二十次全国代表大会上的报告[EB/OL]. (2022-10-25) [2024-09-28]. http://www.news.cn/politics/cpc20/2022-10/25/c_1129079429.htm.
- [2] 中国互联网络信息中心. 第52次《中国互联网络发展状况统计报告》发布[EB/OL]. (2023-08-28) [2024-09-28]. <https://www.cnnic.cn/n4/2023/0828/c199-10830.html>.
- [3] 国家体育总局. 国家体育总局系统深入开展主题教育工作取得实效[EB/OL]. (2023-09-06) [2024-09-28]. <https://www.sport.gov.cn/n20001280/n20745751/c25991100/content.html>.
- [4] 廖列法, 勒孚刚. 基于LDA模型和分类号的专利技术演化研究[J]. 现代情报, 2017, 37(5): 13-18.
- [5] 王晓晨, 关硕, 于文博, 等. 体育赛事网络舆情的传播特征研究——基于2019年女排世界杯的文本情感分析[J]. 成都体育学院学报, 2020, 46(5): 74-81.
- [6] 贾改革. 政民互动中社会诉求主题挖掘和情感分析[D]. 杭州: 浙江大学, 2023.
- [7] 聂思言, 杨江华. 多维视角下新一代人工智能技术的公众感知研究——以微博平台的ChatGPT讨论为例[J]. 情报杂志, 1-10.
- [8] 陈俊宇. 基于文本挖掘的在线评论应用研究[D]. 武汉: 湖北工业大学, 2020.
- [9] 张雷, 谭慧雯, 张璇, 等. 基于LDA模型的高校师德舆情演化及路径传导研究[J]. 情报科学, 2022, 40(3): 144-151.
- [10] 国家体育总局. 中国攀冰公开赛将办 房山将上演冰瀑芭蕾[EB/OL]. (2017-01-02) [2024-09-28]. <https://www.sport.gov.cn/n20001280/n20745751/n20767239/c21736149/content.html>.
- [11] 王晓晨, 关硕, 于文博, 等. 体育赛事网络舆情的传播特征研究——基于2019年女排世界杯的文本情感分析[J]. 成都体育学院学报, 2020, 46(5): 74-81.
- [12] 杨洸. 社会化媒体舆论的极化和共识——以“广州区伯嫖娼”之新浪微博数据为例[J]. 新闻与传播研究, 2016, 23(2): 66-79, 127.
- [13] 王斌, 黄心怡. 平台环境下主流媒体的情感引导效果与传播机制——以抖音号暖新闻为例[J]. 出版广角, 2024(3): 34-41.